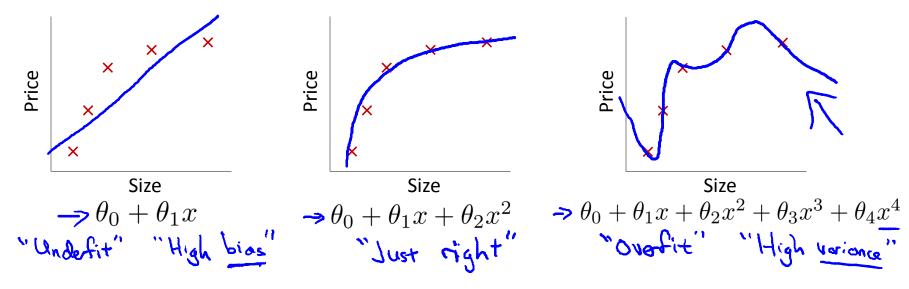


Machine Learning

# Regularization

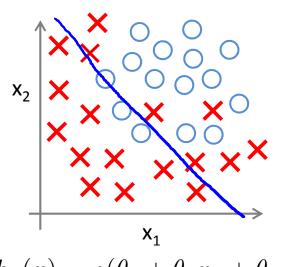
# The problem of overfitting

Example: Linear regression (housing prices)

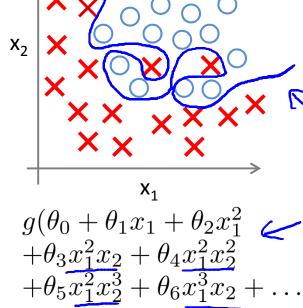


**Overfitting:** If we have too many features, the learned hypothesis may fit the training set very well  $(J(\theta) = \frac{1}{2m} \sum_{i=1}^{m} (h_{\theta}(x^{(i)}) - y^{(i)})^2 \approx 0)$ , but fail to generalize to new examples (predict prices on new examples).

Example: Logistic regression



$$X_{2}$$
 $X_{2}$ 
 $X_{3}$ 
 $X_{4}$ 
 $X_{5}$ 
 $X_{5}$ 
 $X_{5}$ 
 $X_{1}$ 



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

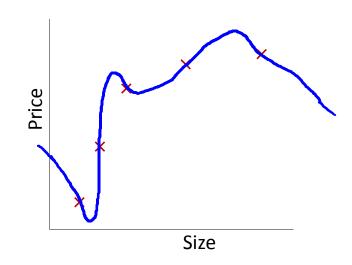
$$(g = \text{sigmoid function})$$

$$g(\theta_{0} + \theta_{1}x_{1} + \theta_{2}x_{2} + \theta_{3}x_{1}^{2} + \theta_{4}x_{2}^{2} + \theta_{5}x_{1}x_{2})$$

$$+\theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 +\theta_5 x_1^2 x_2^3 + \theta_6 x_1^3 x_2 + \dots )$$

#### Addressing overfitting:

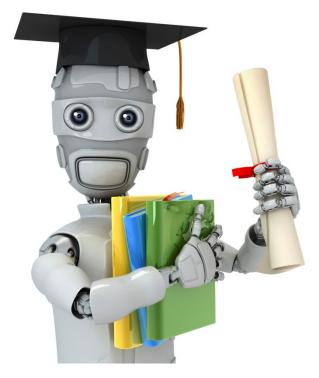
```
x_1 =  size of house
x_2^- no. of bedrooms
x_3 = \text{ no. of floors}
x_4 = age of house
x_5 = average income in neighborhood
x_6 = \text{kitchen size}
x_{100}
```



#### Addressing overfitting:

#### Options:

- 1. Reduce number of features.
- Manually select which features to keep.
- —> Model selection algorithm (later in course).
- 2. Regularization.
  - $\rightarrow$  Keep all the features, but reduce magnitude/values of parameters  $\theta_i$ .
    - Works well when we have a lot of features, each of which contributes a bit to predicting y.



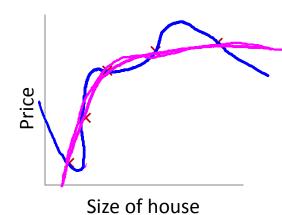
# Regularization

### Cost function

Machine Learning

#### Intuition





$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

Suppose we penalize and make  $\theta_3$ ,  $\theta_4$  really small.

#### Regularization.

Small values for parameters  $\theta_0, \theta_1, \dots, \theta_n \leftarrow$ 

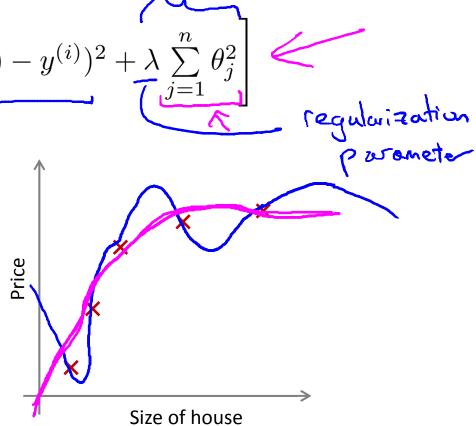
- "Simpler" hypothesis
- Less prone to overfitting <</li>

#### Housing:

- Features:  $x_1, x_2, \dots, x_{100}$
- Parameters:  $\theta_0, \theta_1, \theta_2, \dots, \theta_{100}$

$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^{m} (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \right]$$

#### Regularization.



In regularized linear regression, we choose  $\theta$  to minimize

$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^{m} (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^{n} \theta_j^2 \right]$$

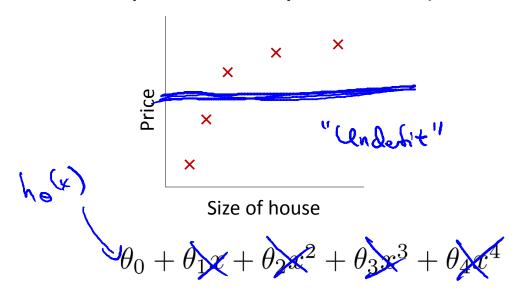
What if  $\lambda$  is set to an extremely large value (perhaps for too large for our problem, say  $\lambda=10^{10}$  )?

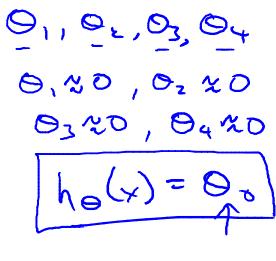
- Algorithm works fine; setting  $\lambda$  to be very large can't hurt it
- Algortihm fails to eliminate overfitting.
- Algorithm results in underfitting. (Fails to fit even training data well).
- Gradient descent will fail to converge.

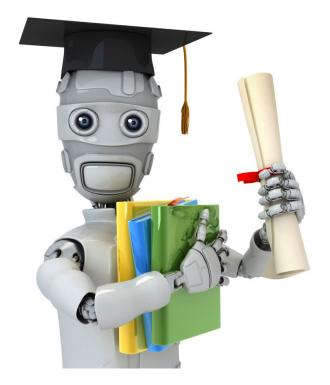
In regularized linear regression, we choose  $\theta$  to minimize

$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^{m} (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^{n} \theta_j^2 \right]$$

What if  $\lambda$  is set to an extremely large value (perhaps for too large for our problem, say  $\lambda=10^{10}$ )?







Machine Learning

# Regularization

Regularized linear regression

#### Regularized linear regression

$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^{m} (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \left( \sum_{j=1}^{n} \theta_j^2 \right) \right]$$

$$\frac{\min_{\theta} J(\theta)}{\uparrow}$$

#### **Gradient descent**



$$\bigcirc$$
,  $\bigcirc$ ,  $\bigcirc$ ,  $\bigcirc$ n

$$= \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^{m}$$

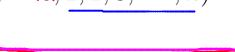
$$\begin{array}{c|c}
 i & \xrightarrow{i=1} & m \\
 \hline
 1 & \xrightarrow{m} & (i) \\
 \end{array}$$

$$\frac{\theta_j}{\tau} := \frac{\theta_j}{\tau}$$

$$\frac{1}{m} \sum_{i=1}^{m} (h_{\theta}(x^{(i)}) - y^{(i)}) x_{j}^{(i)}$$

$$(j = \mathbf{X}, 1, 2, 3, \dots, n)$$







$$= \theta_{j}(1 - \alpha \frac{\lambda}{m}) - \alpha \frac{1}{m} \sum_{i=1}^{m} (h_{\theta}(x^{(i)}) - y^{(i)}) x_{j}^{(i)}$$



#### **Normal equation**

$$X = \begin{bmatrix} (x^{(1)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix}$$

$$\Rightarrow \min_{\theta} J(\theta)$$

$$\Rightarrow 0 = (x^T \times + \lambda)$$

$$\Rightarrow \sum_{\theta \in \mathcal{I}} (x^{(m)})^T$$

$$\Rightarrow \sum_{\theta \in \mathcal{I}} (x^$$

#### Non-invertibility (optional/advanced).

Suppose 
$$m \le n$$
, (#examples) (#features)

$$\theta = \underbrace{(X^T X)^{-1} X^T y}_{\text{Non-invertible / Singular}}$$

If 
$$\frac{\lambda > 0}{\theta} = \left( X^T X + \lambda \begin{bmatrix} 0 & 1 & 1 & 1 \\ & 1 & & \\ & & \ddots & 1 \end{bmatrix} \right)^{-1} X^T y$$

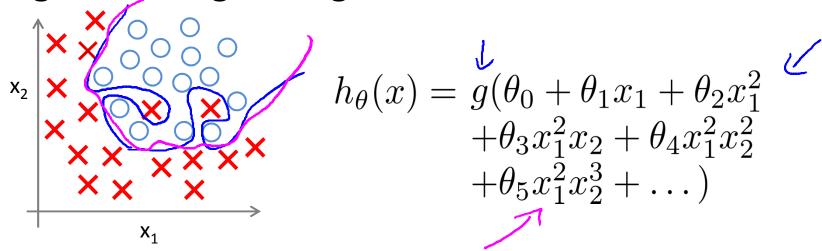


#### Machine Learning

# Regularization

Regularized logistic regression

#### Regularized logistic regression.



#### Cost function:

$$\Rightarrow J(\theta) = -\left[\frac{1}{m} \sum_{i=1}^{m} y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))\right]$$

$$+ \frac{\lambda}{2m} \sum_{j=1}^{n} \mathfrak{S}_{j} \mathfrak{S}_{j}$$
Andrew Andrew

#### **Gradient descent**

Repeat {

$$\theta_{0} := \theta_{0} - \alpha \frac{1}{m} \sum_{i=1}^{m} (h_{\theta}(x^{(i)}) - y^{(i)}) x_{0}^{(i)}$$

$$\theta_{j} := \theta_{j} - \alpha \left[ \frac{1}{m} \sum_{i=1}^{m} (h_{\theta}(x^{(i)}) - y^{(i)}) x_{j}^{(i)} - \frac{\lambda}{m} \Theta_{j} \right]$$

$$\left( \frac{1}{\sqrt{j} = \mathbf{X}, 1, 2, 3, \dots, n} \right)$$

$$\frac{\lambda}{\lambda \Theta_{j}} \mathcal{I}(\Theta)$$

$$h_{\Theta}(\mathbf{x}) = \frac{1}{|\mathbf{x}|^{2} - \mathbf{x}^{2}}$$

#### Advanced optimization

I minunce (e coetendium)? Toot theta(1) <

$$jVal = [code to compute J(\theta)];$$

$$J(\theta) = \left[ \begin{array}{c} \text{code to compute } J(\theta) \\ \end{array} \right];$$

$$J(\theta) = \left[ -\frac{1}{m} \sum_{i=1}^{m} y^{(i)} \log (h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log 1 - h_{\theta}(x^{(i)}) \right] + \left[ \frac{\lambda}{2m} \sum_{j=1}^{n} \theta_{j}^{2} \right]$$

gradient (1) = [code to compute 
$$\frac{\partial}{\partial \theta_0} J(\theta)$$
];

$$\frac{1}{m}\sum_{i=1}^{m}(h_{\theta}(x^{(i)})-y^{(i)})x_{0}^{(i)} \leftarrow$$

$$\Rightarrow \text{gradient (2)} = [\text{code to compute } \frac{\partial}{\partial\theta_{1}}J(\theta)];$$

$$\left(\begin{array}{c|c} \frac{1}{m}\sum_{i=1}^{m}(h_{\theta}(x^{(i)})-y^{(i)})x_{1}^{(i)} - \frac{\lambda}{m}\theta_{1} & \longleftarrow \end{array}\right)$$

$$\rightarrow$$
 gradient (3) = [code to compute  $\frac{\partial}{\partial \theta_2} J(\theta)$ ];

$$\frac{1}{m} \sum_{i=1}^{m} (h_{\theta}(x^{(i)}) - y^{(i)}) x_2^{(i)} - \frac{\lambda}{m} \theta_2$$

gradient (n+1) = [code to compute  $\frac{\partial}{\partial \theta_n} J(\theta)$ ];